

Regression with Gaussian Process Networks

Manel Martínez-Ramón

ECE, UNM

October, 2018

- The solution of the dual parameters is $\boldsymbol{\alpha} = (\mathbf{K} + \sigma_n^2 \mathbf{I})^{-1} \mathbf{y}$ where

$$\mathbf{K} = \boldsymbol{\Phi}^\top \boldsymbol{\Sigma}_p \boldsymbol{\Phi}$$

which defines kernel

$$k(\mathbf{x}, \mathbf{z}) = \boldsymbol{\varphi}(\mathbf{x})^\top \boldsymbol{\Sigma}_p \boldsymbol{\varphi}(\mathbf{z}) = \boldsymbol{\varphi}(\mathbf{x})^\top \boldsymbol{\Sigma}_p^{1/2} \boldsymbol{\Sigma}_p^{1/2} \boldsymbol{\varphi}(\mathbf{z})$$

which is a dot product of vectors linearly transformed by matrix $\boldsymbol{\Sigma}_p^{1/2}$.

- The choice of this matrix is implicit in the choice of the kernel.

- With the previous definition of the kernel, the regression can be written as

$$\begin{aligned}\bar{f}_* &= \boldsymbol{\varphi}(x^*)^\top \bar{\mathbf{w}}' \\ &= \boldsymbol{\varphi}(x^*)^\top \boldsymbol{\Sigma}_p \boldsymbol{\Phi} \boldsymbol{\alpha} \\ &= \mathbf{k}^\top(\mathbf{x}^*) \boldsymbol{\alpha} \\ &= \mathbf{k}^\top(\mathbf{x}^*) (\mathbf{K} + \sigma_n^2 \mathbf{I})^{-1} \mathbf{y}\end{aligned}$$

where $\mathbf{k}^\top(\mathbf{x}^*) = \{k(\mathbf{x}^*, \mathbf{x}[1]) \cdots k(\mathbf{x}^*, \mathbf{x}[N])\}$ is the column vector of dot products between the test data \mathbf{x}^* and all training data $\mathbf{x}[n]$ into the Hilbert space.

- Expression

$$\bar{f}_* = \mathbf{k}^\top(\mathbf{x}^*) (\mathbf{K} + \sigma_n^2 \mathbf{I})^{-1} \mathbf{y}$$

is the *expectation of the Gaussian process* at sample \mathbf{x}^* conditional to \mathbf{X}, \mathbf{y} .

- We can compute the variance of the process at this point simply using the obtained definition of the kernel. Indeed:

$$\sigma_{f_*}^2 = \varphi(\mathbf{x})^{*\top} \mathbf{A}^{-1} \varphi(\mathbf{x}^*)$$

with $\mathbf{A} = \sigma_n^{-2} \boldsymbol{\Phi} \boldsymbol{\Phi}^\top + \boldsymbol{\Sigma}_p^{-1}$.

- We compute the inverse of \mathbf{A} using the matrix inversion lemma:

$$(\mathbf{U}\mathbf{W}\mathbf{V} + \mathbf{Z})^{-1} = \mathbf{Z}^{-1} - \mathbf{Z}^{-1}\mathbf{U}(\mathbf{V}^\top\mathbf{Z}^{-1}\mathbf{U} + \mathbf{W}^{-1})^{-1}\mathbf{V}^\top\mathbf{Z}^{-1}$$

Since $\mathbf{A} = \sigma_n^{-2}\boldsymbol{\Phi}\boldsymbol{\Phi}^\top + \boldsymbol{\Sigma}_p^{-1}$ we can identify

$$\mathbf{Z} = \boldsymbol{\Sigma}_p^{-1}, \quad \mathbf{U} = \mathbf{V} = \boldsymbol{\Phi}, \quad \mathbf{W} = \sigma_m^{-2}\mathbf{I}$$

from which

$$\mathbf{A}^{-1} = \boldsymbol{\Sigma}_p - \boldsymbol{\Sigma}_p\boldsymbol{\Phi}(\mathbf{K} + \sigma_n^2\mathbf{I})^{-1}\boldsymbol{\Phi}^\top\boldsymbol{\Sigma}_p$$

- Now, from expression

$$\sigma_{f_*}^2 = \varphi(\mathbf{x})^{*\top} \mathbf{A}^{-1} \varphi(\mathbf{x}^*)$$

we obtain

$$\begin{aligned} \sigma_{f_*}^2 &= \varphi(\mathbf{x})^{*\top} \left(\Sigma_p - \Sigma_p \Phi (\mathbf{K} + \sigma_n^2 \mathbf{I})^{-1} \Phi^\top \Sigma_p \right) \varphi(\mathbf{x}^*) \\ &= \varphi(\mathbf{x})^{*\top} \Sigma_p \varphi(\mathbf{x}^*) - \varphi(\mathbf{x})^{*\top} \Sigma_p \Phi (\mathbf{K} + \sigma_n^2 \mathbf{I})^{-1} \Phi^\top \Sigma_p \varphi(\mathbf{x}^*) \end{aligned}$$

- Using the definition of the kernel, the variance is

$$\sigma_{f_*}^2 = k(\mathbf{x}^*, \mathbf{x}^*) - \mathbf{k}(\mathbf{x}^*)^\top (\mathbf{K} + \sigma_n^2 \mathbf{I})^{-1} \mathbf{k}(\mathbf{x}^*)$$

- The predictive mean and variance of the Gaussian process in a Hilbert space defined by kernel $k(\mathbf{x}, \mathbf{x})$ evaluated at \mathbf{x}^* are

$$\bar{f}_* = \mathbf{k}(\mathbf{x}^*) (\mathbf{K} + \sigma_n^2 \mathbf{I})^{-1} \mathbf{y}$$

and

$$\sigma_{f_*}^2 = k(\mathbf{x}^*, \mathbf{x}^*) - \mathbf{k}(\mathbf{x}^*)^\top (\mathbf{K} + \sigma_n^2 \mathbf{I})^{-1} \mathbf{k}(\mathbf{x}^*)$$

where $\mathbf{k}(\mathbf{x}^*) = \{k(\mathbf{x}^*, \mathbf{x}[1]) \cdots k(\mathbf{x}^*, \mathbf{x}[N])\}^\top$.

Proposition: *The kernel matrix \mathbf{K} is the covariance matrix of the Gaussian process estimation of $f(\mathbf{x}[n])$.*

Proof: We assume that the process has zero mean, this is $\mathbb{E}[f(\mathbf{x}[n])] = 0$. Then, straightforwardly, the covariance matrix is

$$\begin{aligned}\mathbb{E}[f(\mathbf{x}[n])f(\mathbf{x}[m])] &= \mathbb{E}[\boldsymbol{\varphi}(\mathbf{x}[n])^\top \mathbf{w}' \mathbf{w}'^\top \boldsymbol{\varphi}(\mathbf{x}[m])] \\ &= \boldsymbol{\varphi}^\top(\mathbf{x}[n]) \mathbb{E}[\mathbf{w}' \mathbf{w}'^\top] \boldsymbol{\varphi}(\mathbf{x}[m]) \\ &= \boldsymbol{\varphi}^\top(\mathbf{x}[n]) \boldsymbol{\Sigma}_p \boldsymbol{\varphi}(\mathbf{x}[m]) = k(\mathbf{x}[n], \mathbf{x}[m])\end{aligned}$$

Corollary There is an error $\varepsilon[n]$ in the estimation

$$y[n] = f(\mathbf{x}[n]) + \varepsilon[n]$$

which is modelled as AWGN, independent of $\mathbf{x}[n]$ and with variance σ_n^2 . Then, the covariance of the regressors is

$$\mathbb{E}[y[n]y[m]] = k(\mathbf{x}[n], \mathbf{x}[m]) + \sigma_n^2 \delta(m - n)$$

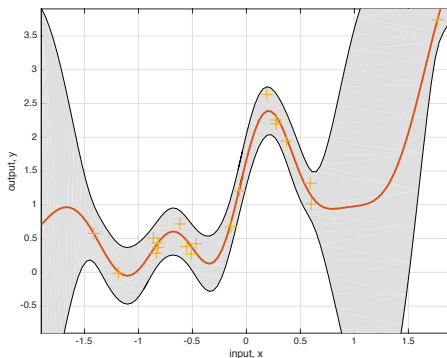
so the covariance matrix of the process is

$$\mathbf{C}_{yy} = \mathbf{K} + \sigma_n^2 \mathbf{I}$$

- The previous result gives an alternative interpretation of the predictive likelihood with the mean and variance summarized in slide 11.
- Assuming a set of training data \mathbf{X} and corresponding regressors \mathbf{y} , and a test data \mathbf{x}^* , the joint process \mathbf{y}, f_* is, from the previous results

$$\begin{bmatrix} \mathbf{y} \\ f_* \end{bmatrix} \sim \mathcal{N} \left(\mathbf{0}, \begin{bmatrix} \mathbf{K} + \sigma_n^2 \mathbf{I} & \mathbf{k}(\mathbf{x}^*) \\ \mathbf{k}^\top(\mathbf{x}^*) & k(\mathbf{x}^*, \mathbf{x}^*) \end{bmatrix} \right)$$

Using the Bayes rule, one can compute the pdf of f_* conditional to \mathbf{X}, \mathbf{y} , which is a Gaussian with the mean and variance of slide 7.



Example included in the software provided in Rasmussen et al, 2006. A set of samples is generated from a filtered Gaussian distribution. The line represents the predictive mean. The band represents the standard deviation of the prediction.

After this video, students should be able to:

- Identify the mean and variance of the predictive posterior in a Kernel Gaussian Process
- Reproduce the proofs for both expressions.
- Prove that the kernel matrix is equal to the covariance matrix of the training regressors.